# GigaScience

## Genome Annotation Generator: A simple tool for generating and correcting WGS annotation tables for NCBI submission

### --Manuscript Draft--

| Manuscript Number: | GIGA-D-17-00030R2 |
|---|---|
| Full Title: | Genome Annotation Generator: A simple tool for generating and correcting WGS annotation tables for NCBI submission |
| Article Type: | Technical Note |
| Funding Information: | Agricultural Research Service | Dr Scott M Geib |

| Abstract: | One of the most overlooked, yet critical components of a whole genome sequencing project is the submission and curation of the data to a genomic repository, most commonly NCBI. While large genome centers or genome groups have developed software tools for post-annotation assembly filtering, annotation, and conversion into NCBI's annotation table format, these tools typically require back-end setup and connection to an SQL database and/or some knowledge of programming (Perl, Python) to implement. With whole genome sequencing becoming commonplace, genome sequencing projects are moving away from the genome centers, and into the ecology or biology lab, where much less resources are present to support the process of genome assembly curation. To fill this gap, we developed software to assess, filter, transfer annotations, and convert a draft genome assembly and annotation set into NCBI annotation table (.tbl) format, facilitating submission to NCBI Genome Assembly database. This software has no dependencies, is compatible across platforms, and utilizes a simple command line to perform a variety of simple and complex post-analysis, pre-NCBI submission WGS project tasks. |
|---|---|

| Corresponding Author: | Scott M Geib<br><br>UNITED STATES |
|---|---|
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Scott M Geib |
| First Author Secondary Information: | |
| Order of Authors: | Scott M Geib |
| | Brian Hall |
| | Theodore DeRego |
| | Forest T Bremer |
| | Kyle Cannoles |
| | Sheina B Sim |
| Order of Authors Secondary Information: | |

| Response to Reviewers: | Hi, I have modified the manuscript to reference a GigaDB entry with a stable release for software access rather than just referencing a github page |
|---|---|
| Additional Information: | |
| Question | Response |
| Are you submitting this manuscript to a special series or article collection? | No |

| | |
|---|---|
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our [Minimum Standards Reporting Checklist](). Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite [Research Resource Identifiers]() (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our [Minimum Standards Reporting Checklist]()? | Yes |
| **Availability of data and materials**<br><br>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories]() (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.<br><br>Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist]()? | Yes |

# Genome Annotation Generator: A simple tool for generating and correcting WGS annotation tables for NCBI submission

Scott M. Geib [1][*][†], Brian Hall [2][†], , Theodore Derego [1], Forest T. Bremer,[2], Kyle Cannoles[1,3], and Sheina B. Sim [1]

[1]Tropical Plant Protection Research Unit, USDA-ARS Daniel K. Inouye U.S. Pacific Basin Agricultural Research Center, Hilo, HI, 96720, USA

[2]Plant and Environmental Protection Science, University of Hawaii at Manoa, Honolulu, HI, 96822, USA.

[3]Department of Computer Science and Engineering, University of Hawaii at Hilo, Hilo, HI, 96720, USA.

*Corresponding author [†]Authors contributed equally

email: SMG: scott.geib@ars.usda.gov BH:bhall7@hawaii.edu TD:t.derego@yahoo.com FTB:forestb@hawaii.edu KC:kylecann@hawaii.edu SBS:sheina.sim@ars.usda.gov

February 21, 2018

## Abstract

**Background**

One of the most overlooked, yet critical components of a whole genome sequencing project is the submission and curation of the data to a genomic repository, most commonly NCBI. While large genome centers or genome groups have developed software tools for post-annotation assembly filtering, annotation, and conversion into NCBI's annotation table format, these tools typically require back-end setup and connection to an SQL database and/or some knowledge of programming (Perl, Python) to implement. With whole genome sequencing becoming commonplace, genome sequencing projects are moving away from the genome centers, and into the ecology or biology lab, where much less resources are present to support the process of genome assembly curation. To fill this gap, we developed software to assess, filter, transfer annotations, and convert a draft genome assembly and annotation set into NCBI annotation table (.tbl) format, facilitating submission to NCBI Genome Assembly database. This software has no dependencies, is compatible across platforms, and utilizes a simple command to perform a variety of simple and complex post-analysis, pre-NCBI submission WGS project tasks.

**Findings**

The Genome Annotation Generator is a consistent and user-friendly bioinformatics tool that can be used to generate a *.tbl* file that is consistent with the NCBI submission pipeline.

**Conclusions**

The Genome Annotation Generator achieves the goal of providing a publicly available tool that will facilitate the submission of annotated genome assemblies to NCBI. It is useful for any individual researcher or research group who wishes to submit a genome assembly of their study system to NCBI.

**Keywords:** Genome curation; annotation; and whole-genome sequencing project

# 1 Introduction

While ever-improving sequencing technology and assembly software enable the collection of raw sequences for genome assembly and structural annotation, further steps need to be taken to ensure the quality and completeness of a WGS project for submission to the National Center for Biotechnology Information (NCBI) or other data repositories [36]. To submit a genome to the NCBI for curation, it must be converted to the NCBI annotation table format (*.tbl*). With a genome assembly project consisting of thousands of sequences demarcated by hundreds of thousands of structural annotations, this task clearly requires automation. However, there is currently no freely available tool which performs rapid and controlled conversion of a genome assembly and associated structural annotations into a *.tbl* format in addition to allowing for editing, modification, and revision of the content of the project. Moreover, the typical assembly and draft annotation contains some degree of questionable or erroneous data which requires correction or omission. It may also be desirable to add functional annotations to the submission and integrate results from InterProScan, BLAST homology to curated databases, or ontology terms generated by other tools [20, 5, 22].

The traditional approach used to address these problems is to use Linux command line tools or write custom scripts which modify and filter the genome using a scripting language such as Perl or Python [4, 32, 15] or large scale genomic database systems [25]. This method may not be easily or readily reproducible, or it may be entirely beyond the ability of an investigator who has less familiarity with generating custom scripts *de novo*. Even amongst those researchers who use best practices to write clean, well-tested, and reusable scripts to accomplish these tasks, doing requires a large amount of duplicated effort. For this reason, the Genome Annotation Generator (GAG) was written to provide a straightforward and consistent tool for addressing the most common errors in genome assemblies, adding functional annotations from disparate sources, and producing an NCBI submission-ready annotation *.tbl* file. In addition, the software provides a means for integrating existing functional annotations and marking annotations that require manual curation or review. All of these tasks are done through an intuitive command line program requiring only basic unix skills, and has no required dependencies or packages. The program GAG facilitates the submission of whole genome sequencing (WGS) projects to NCBI as well as provide a standardized utility and workflow that fosters consistency

2

between projects. Due to emerging genome sequencing initiatives such as the 5,000 Insect Genomes Initiative (i5K), the Plant Genome Initiative, and Genome 10K [19, 26], many independent research groups which are not specialized in genome annotation and analysis are generating large genomic datasets and performing genome sequencing projects within their lab. This program can assist in ensuring quality and consistency of data for new genome biologists.

# 2   Overview

The GAG program is a command line Python program, written in Python 2.7 and requiring no additional outside programs or packages to run. The user directs the program to the genome *.fasta* file and a *.gff3* file containing structure annotations. In addition, a number of options can be used to fix possible errors, flag or remove features (i.e. genomic elements described in the *.gff* structural anntoation file) based on selected criteria, add functional annotations, trim regions of the genome out of the assembly, and, of course, write the genome to NCBI *.tbl* file format. In addition, changes made to the genome annotation, functional annotations added, or flags requesting manual review are also annotated back to the *.gff3* structural annotation file, and the original fasta file is corrected as needed. When the user issues commands to modify the genome, e.g. to remove short introns, the statistics will display two columns, representing the original and modified genomes. This allows for stepwise and documented filtering and review to occur, and interactions between GAG and visual genome review tools (e.g. Artemis, Apollo, GBrowse, JBrowse etc) [35, 29, 21, 31, 28].

# 3   Methods

As an example, we consider a possible work-flow for a user wishing to prepare a genome for submission to the NCBI Eukaryotic WGS Database. The user has a scaffolded genome assembly produced by one of many whole genome assemblers [2, 16, 30] in *.fasta* file format and a corresponding GFF3 feature file [9, 8] containing structural annotations resulting from an automated annotation pipeline or predictors such as Maker, Evidence Modeler, Jigsaw, or others e.g. [3, 18, 1, 17, 33, 34, 7, 10]. The approach would be to first possibly generate functional annotations of predicted genes if this is desired, using whatever approach the user is interested in, and then using the genome and annotations with GAG. After using GAG to remove or flag features of interest, the user then may then further investigate flagged features in a genome browser by loading the output of GAG, edit, and then perform further filtering in GAG, and iterate through this process until a final draft genome product is generated. Finally GAG writes a NCBI table file, on which *tbl2asn* is run for submission to NCBI. This may identify regions of the genome

3

that need to be trimmed, due to possible adapter contamination in the genome, or low quality sequence. Any errors generated by *tbl2asn* can then be corrected in GAG, the genome trimmed, until an error free submission is generated.

To use GAG, the user creates a folder containing the genome files (or links to them) and runs *gag.py* from the terminal, with the *.fasta* and *.gff3* files. GAG will write a statistics file, containing infomraiton on the number of each feature type, lengths, and other information that may be useful for the submitter. In our experience, automated genome annotation software frequently produces assemblies containing introns as short as 1 base pair long; if any such features are present, GAG can be run to detect them. It is important to note that while NCBI requires short introns to be removed, cutoffs recommended by NCBI may be more stringent that what you want, as they are set to reduce erroneous data being entered into NCBI. For example, prediction of single base introns might not be errors, and represent true data. It is up to the user to dictate what cutoffs they want to define to remove or flag for manual review. To address these short introns, the user simply applies option *-ris* (Remove_Intron_Shorter_Than) with a value of *10*. GAG will discard any mRNA containing an intron shorter than the minimum of ten. A comparison of the genome content before and after removal is printed to the *.stats* file. If the user instead wishes to only flag features that meet these criteria and not remove them, alternatively the *-fis* (Flag_Introns_Shorter_Than) option could be used, which instead adds a *GAG_FLAG* annotation to the attributes column of the *.gff3* file describing the reason for flagging, allowing manual review of flagged features in a genome browser. GAG will automatically update all parent and child features (gene or CDS entries) to reflect removal of mRNA features. A list of available flag or removal options are listed in Table 1.

Another review for submission might be that all coding regions be a minimum length. For this example we use 150 base pairs in length, which is suggested by NCBI [11, 12]. To add this additional level of filtering, a second option can be used: *-rcs 150*, to Remove_CDS_Shorter_Than 150 bp. When the genome is written to the output folder, GAG will write a file called *genome.removed.gff* containing all the features left out of the final version). It is important to remember that CDS cutoffs at 150 bp will possibly remove some biologically correct amino acids.

GAG supports two straightforward correction, or fix tools. If the user's GFF3 file does not explicitly indicate the presence of start and stop codons, or if there is reason to believe there are errors in ORF prediction in the provided GFF file, GAG can add start and stop feature to the GFF file. The user simply issues the command with the option *–fix_start_stop* and these features will be added to the GFF3 file, and their existence noted in the table file. A second issue that can arise in a draft genome assembly is for a contig or scaffold to have a string of ambiguous bases (N's) at the very beginning or end of the contig. These should be removed from the assembly, and can be using the –

*fix_terminal_ns* option, as they can be mis-interpreted as scaffold gaps. Removing these regions from the genome though, will disrupt the parity between coordinates in the *.fasta* genome file and the *.gff3* annotation file. GAG will automatically update coordinates in the *.gff3* file to reflect any regions removed from the sequence file. During execution of *tbl2asn* or submission to NCBI, it may be identified that regions of the genome may be contaminated with microbial, vector, or sequencing adapter sequence as part of the "contaminate screen" step. A *.bed* formatted file can be supplied with the *-trim* option, containing regions of the assembly to exclude, either ranges within a contig or scaffold, or an entire scaffold. GAG will update both the *.fasta* and *.gff3* files so that coordinate are still synchronized. This is a particularly difficult operation to perform without a specialized tool.

At present, GAG has simple commands to remove or flag introns, exons, coding regions and genes based on minimum or maximum lengths, which will also edit or remove any parent or child feature from the annotation file so as not to create incomplete feature annotations. It can also remove features from a list, which is useful for cases where a genome submission is rejected and a list of invalid mRNAs and genes provided. In addition, all discarded features are retained in a "genome.removed.gff" file and the entire editing session is documented so that the user can retain the filtering criteria used on the particular dataset.

GAG supports two methods to add functional annotations to a genome. First, it can read an annotated GFF3 file containing gene names, protein products, cross-references to databases, and ontology terms following GFF3 qualified nomenclature in the *attribute* column of the GFF3 file [24, 23, 27, 6]. Any annotations present will be automatically carried over to the NCBI feature table file. For users with annotations from another source, GAG can read them from a simple tab-delimited file. The annotations supported by the current version of GAG are *Name* (for genes), *Dbxref*, *Ontology_term* and *product* (for descriptive mRNA products). These are also written to a new GFF3 file, so GAG can be utilized as a tool to also functionally annotate a GFF3 file. Detailed instructions for running GAG, examples for each of the main functions (e.g. removing features, adding start and stop codons, trimming features, adding annotations) as well as formats and conversion tools for functional annotations are available on the GAG software website webpage: `http://genomeannotation.github.io/GAG/` and a stable release (version 2.0.1) is available in an accompanying GigaScience Database entry [13, 14] and is shared through SciCrunch under RRID SCR_016053.

# 4 Implementation

GAG is written in Python 2.7. It has no dependencies beyond the standard library. The program is modular, abstracting biological concepts such as Sequence, Gene and CDS

5

into classes which may be incorporated into other software tools. In addition, the code is covered by a suite of unit and integration tests, allowing developers to modify or add to the code base with reduced risk of introducing errors. It should be easily executable by the novice programmer with only basic command-line experience, but also powerful enough to be implemented within robust genomic data processing pipelines.

# 5   Conclusion

GAG can be easily expanded in the future to support more specific needs of researchers, less common annotation types, and integrate conversion of common functional annotation output formats (e.g. InterProScan, BLAST, Blast2Go) for addition to NCBI annotation table formats. Currently, GAG is an intermediate, but critical tool, between a simple format conversion tool and more sophisticated annotation editors. In future developments of GAG, we plan to allow the integration of multiple lines of evidence supporting gene models to help users discriminate apparently high quality annotations from annotations with little support or possible errors. This could rapidly improve and standardize manual annotation efforts in systems and user groups that are not integrated into genome center annotation pipelines.

# 6   Data and Software

Availability and requirements

- Project name: Genome Annotation Generator

- Project home page: `https://github.com/genomeannotation/GAG`

- Operating systems: Linux, Windows, OS

- Programming language: Python

- Other requirements: Python 2

- License: MIT

- SciCrunch RRID: SCR_016053

Availability of Supporting Data Snapshots of the software and accompanying files are available from the GigaScience GigaDB database, and the algorithm is also demonstrated in Code Ocean repository `https://codeocean.com/2018/02/06/genome-annotation-generator-ncbi-for-submission/` [13] [14].

# 7 Declarations

## 7.1 Competing Interests

The authors declare that they have no competing interests

## 7.2 Funding

## 7.3 Authors' contributions

SMG conceived software concept. BH, TD, and SMG designed and wrote software. BH, SMG, and SBS wrote manuscript.

## 7.4 Acknowledgements

Table 1: **Options for GAG**

| Option | Type of function | Description |
| --- | --- | --- |
| *-a* <annotation file> | Annotate | Adds functional annotations present in annotation file to *.gff* and *.tbl* |
| *-t* <*.bed* file> | Trim | Removes regions of genome indicated in *.bed* file from *.fasta* and *.gff3* |
| *-fix_start_stop* <no value> | Fix | Adds or corrects start and stop codon features to *.gff3* |
| *-fix_terminal_ns* <no value> | Fix | Removes any trailing ends from contig ends in assembly, updates *.gff3* coordinates |
| *-rcs* <integer> | Remove | Remove CDS shorter than <integer> |
| *-rcl* <integer> | Remove | Remove CDS longer than <integer> |
| *-res* <integer> | Remove | Remove exons shorter than <integer> |
| *-rel* <integer> | Remove | Remove exons longer than <integer> |
| *-ris* <integer> | Remove | Remove introns shorter than <integer> |
| *-ril* <integer> | Remove | Remove introns longer than <integer> |
| *-rgs* <integer> | Remove | Remove genes shorter than <integer> |
| *-rgl* <integer> | Remove | Remove genes longer than <integer> |
| *-fcs* <integer> | Flag | Flag CDS shorter than <integer> |
| *-fcl* <integer> | Flag | Flag CDS longer than <integer> |
| *-fes* <integer> | Flag | Flag exons shorter than <integer> |
| *-fel* <integer> | Flag | Flag exons longer than <integer> |
| *-fis* <integer> | Flag | Flag introns shorter than <integer> |
| *-fil* <integer> | Flag | Flag introns longer than <integer> |
| *-fgs* <integer> | Flag | Flag genes shorter than <integer> |
| *-fgl* <integer> | Flag | Flag genes longer than <integer> |

# References

[1] Jonathan E. Allen and Steven L. Salzberg. Jigsaw: integration of multiple sources of evidence for gene prediction. *Bioinformatics*, 21(18):3596–3603, 2005.

[2] J. Butler, I. MacCallum, M. Kleber, I. A. Shlyakhter, M. K. Belmonte, E. S. Lander, C. Nusbaum, and D. B. Jaffe. Allpaths: de novo assembly of whole-genome shotgun microreads. *Genome Res*, 18(5):810–20, 2008.

[3] Brandi L. Cantarel, Ian Korf, Sofia M. C. Robb, Genis Parra, Eric Ross, Barry Moore, Carson Holt, Alejandro Sanchez Alvarado, and Mark Yandell. Maker: An easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Research*, 18(1):188–196, 2008.

[4] Peter J. A. Cock, Tiago Antao, Jeffrey T. Chang, Brad A. Chapman, Cymon J. Cox, Andrew Dalke, Iddo Friedberg, Thomas Hamelryck, Frank Kauff, Bartek Wilczynski, and Michiel J. L. de Hoon. Biopython: freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423, 2009.

[5] Ana Conesa, Stefan Götz, Juan Miguel García-Gómez, Javier Terol, Manuel Talón, and Montserrat Robles. Blast2go: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics*, 21(18):3674–3676, 2005.

[6] The Gene Ontology Consortium. Expansion of the gene ontology knowledgebase and resources. *Nucleic Acids Research*, 45(D1):D331–D338, 2017.

[7] Val Curwen, Eduardo Eyras, T. Daniel Andrews, Laura Clarke, Emmanuel Mongin, Steven M.J. Searle, and Michele Clamp. The ensembl automatic gene annotation system. *Genome Research*, 14(5):942–950, 2004.

[8] K. Eilbeck and S. E. Lewis. Sequence ontology annotation guide. *Comp Funct Genomics*, 5(8):642–7, 2004.

[9] Karen Eilbeck, Suzanna Lewis, Christopher Mungall, Mark Yandell, Lincoln Stein, Richard Durbin, and Michael Ashburner. The sequence ontology: a tool for the unification of genome annotations. *Genome Biology*, 6(5):R44, 2005.

[10] Christine G. Elsik, Kim C. Worley, Anna K. Bennett, Martin Beye, Francisco Camara, Christopher P. Childers, Dirk C. de Graaf, Griet Debyser, Jixin Deng, Bart Devreese, Eran Elhaik, Jay D. Evans, Leonard J. Foster, Dan Graur, Roderic Guigo, Katharina Jasmin Hoff, Michael E. Holder, Matthew E. Hudson, Greg J. Hunt, Huaiyang Jiang, Vandita Joshi, Radhika S. Khetani, Peter Kosarev, Christie L.

Kovar, Jian Ma, Ryszard Maleszka, Robin F. A. Moritz, Monica C. Munoz-Torres, Terence D. Murphy, Donna M. Muzny, Irene F. Newsham, Justin T. Reese, Hugh M. Robertson, Gene E. Robinson, Olav Rueppell, Victor Solovyev, Mario Stanke, Eckart Stolle, Jennifer M. Tsuruda, Matthias Van Vaerenbergh, Robert M. Waterhouse, Daniel B. Weaver, Charles W. Whitfield, Yuanqing Wu, Evgeny M. Zdobnov, Lan Zhang, Dianhui Zhu, and Richard A. Gibbs. Finding the missing honey bee genes: lessons learned from a genome upgrade. *BMC Genomics*, 15(1):86, 2014.

[11] National Center for Biotechnology Information. The genbank submissions handbook [internet], 2011.

[12] National Center for Biotechnology Information. Common discrepancy reports, January 2013.

[13] S. M. Geib, B. Hall, T. Derego, F. T. Bremer, K. Cannoles, and S. B. Sim. Supporting data for "genome annotation generator: A simple tool for generating and correcting wgs annotation tables for ncbi submission". *GigaScience Database*, 2017. `http://dx.doi.org/10.5524/100308`.

[14] S. M. Geib, B. Hall, T. Derego, F. T. Bremer, K. Cannoles, and S. B. Sim. Genome annotation generator ncbi for submission [source code]. *CodeOcean*, 2018. `https://doi.org/10.24433/CO.fceb0521-a26d-441f-9fe0-bccc6a250fc9`.

[15] Robert Gentleman, Vincent Carey, Douglas Bates, Ben Bolstad, Marcel Dettling, Sandrine Dudoit, Byron Ellis, Laurent Gautier, Yongchao Ge, Jeff Gentry, Kurt Hornik, Torsten Hothorn, Wolfgang Huber, Stefano Iacus, Rafael Irizarry, Friedrich Leisch, Cheng Li, Martin Maechler, Anthony Rossini, Gunther Sawitzki, Colin Smith, Gordon Smyth, Luke Tierney, Jean Yang, and Jianhua Zhang. Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology*, 5(10):R80, 2004.

[16] Sante Gnerre, Iain MacCallum, Dariusz Przybylski, Filipe J. Ribeiro, Joshua N. Burton, Bruce J. Walker, Ted Sharpe, Giles Hall, Terrance P. Shea, Sean Sykes, Aaron M. Berlin, Daniel Aird, Maura Costello, Riza Daza, Louise Williams, Robert Nicol, Andreas Gnirke, Chad Nusbaum, Eric S. Lander, and David B. Jaffe. High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proceedings of the National Academy of Sciences*, 108(4):1513–1518, 2011.

[17] Brian Haas, Steven Salzberg, Wei Zhu, Mihaela Pertea, Jonathan Allen, Joshua Orvis, Owen White, C Robin Buell, and Jennifer Wortman. Automated eukaryotic gene structure annotation using evidencemodeler and the program to assemble spliced alignments. *Genome Biology*, 9(1):R7, 2008.

[18] Carson Holt and Mark Yandell. Maker2: an annotation pipeline and genome-database management tool for second-generation genome projects. *Bmc Bioinformatics*, 12, 2011.

[19] i5K Consortium. The i5k initiative: Advancing arthropod genomics for knowledge, human health, agriculture, and the environment. *Journal of Heredity*, 104(5):595–600, 2013.

[20] Philip Jones, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish McWilliam, John Maslen, Alex Mitchell, Gift Nuka, Sebastien Pesseat, Antony F. Quinn, Amaia Sangrador-Vegas, Maxim Scheremetjew, Siew-Yit Yong, Rodrigo Lopez, and Sarah Hunter. Interproscan 5: genome-scale protein function classification. *Bioinformatics*, 30(9):1236–1240, 2014.

[21] Eduardo Lee, Gregg A. Helt, Justin T. Reese, Monica C. Munoz-Torres, Chris P. Childers, Robert M. Buels, Lincoln Stein, Ian H. Holmes, Christine G. Elsik, and Suzanna E. Lewis. Web apollo: a web-based genomic annotation editing platform. *Genome Biology*, 14(8):R93, 2013.

[22] Michele Magrane and UniProt Consortium. Uniprot knowledgebase: a hub of integrated protein data. *Database*, 2011, 2011.

[23] Barry Moore, Guozhen Fan, and Karen Eilbeck. Soba: sequence ontology bioinformatics analysis. *Nucleic Acids Research*, 38(suppl 2):W161–W164, 2010.

[24] Christopher J. Mungall, Colin Batchelor, and Karen Eilbeck. Evolution of the sequence ontology terms and relationships. *Journal of Biomedical Informatics*, 44(1):87–93, 2011.

[25] Christopher J. Mungall, David B. Emmert, and The FlyBase Consortium. A chado case study: an ontology-based modular schema for representing genome-associated biological information. *Bioinformatics*, 23(13):i337–i346, 2007.

[26] Genome 10K Community of Scientists. Genome 10k: A proposal to obtain whole-genome sequence for 10,000 vertebrate species. *Journal of Heredity*, 100(6):659–674, 2009.

[27] Martin Reese, Barry Moore, Colin Batchelor, Fidel Salas, Fiona Cunningham, Gabor Marth, Lincoln Stein, Paul Flicek, Mark Yandell, and Karen Eilbeck. A standard variation file format for human genome sequences. *Genome Biology*, 11(8):R88, 2010.

[28] J. T. Robinson, H. Thorvaldsdottir, W. Winckler, M. Guttman, E. S. Lander, G. Getz, and J. P. Mesirov. Integrative genomics viewer. *Nat Biotechnol*, 29(1):24–6, 2011.

[29] K. Rutherford, J. Parkhill, J. Crook, T. Horsnell, P. Rice, M. A. Rajandream, and B. Barrell. Artemis: sequence visualization and annotation. *Bioinformatics*, 16(10):944–5, 2000.

[30] Jared T. Simpson, Kim Wong, Shaun D. Jackman, Jacqueline E. Schein, Steven J.M. Jones, and Inanc Birol. Abyss: A parallel assembler for short read sequence data. *Genome Research*, 19(6):1117–1123, 2009.

[31] Mitchell E. Skinner, Andrew V. Uzilov, Lincoln D. Stein, Christopher J. Mungall, and Ian H. Holmes. Jbrowse: A next-generation genome browser. *Genome Research*, 19(9):1630–1638, 2009.

[32] Jason E. Stajich, David Block, Kris Boulez, Steven E. Brenner, Stephen A. Chervitz, Chris Dagdigian, Georg Fuellen, James G.R. Gilbert, Ian Korf, Hilmar Lapp, Heikki Lehväslaiho, Chad Matsalla, Chris J. Mungall, Brian I. Osborne, Matthew R. Pocock, Peter Schattner, Martin Senger, Lincoln D. Stein, Elia Stupka, Mark D. Wilkinson, and Ewan Birney. The bioperl toolkit: Perl modules for the life sciences. *Genome Research*, 12(10):1611–1618, 2002.

[33] Mario Stanke, Oliver Schoffmann, Burkhard Morgenstern, and Stephan Waack. Gene prediction in eukaryotes with a generalized hidden markov model that uses hints from external sources. *BMC Bioinformatics*, 7(1):62, 2006.

[34] Mario Stanke and Stephan Waack. Gene prediction with a hidden markov model and a new intron submodel. *Bioinformatics*, 19(suppl 2):ii215–ii225, 2003.

[35] Lincoln D. Stein, Christopher Mungall, ShengQiang Shu, Michael Caudy, Marco Mangone, Allen Day, Elizabeth Nickerson, Jason E. Stajich, Todd W. Harris, Adrian Arva, and Suzanna Lewis. The generic genome browser: A building block for a model organism system database. *Genome Research*, 12(10):1599–1610, 2002.

[36] Mark Yandell and Daniel Ence. A beginner's guide to eukaryotic genome annotation. *Nature Reviews Genetics*, 13(5):329–342, 2012.

# Genome Annotation Generator: A simple tool for generating and correcting WGS annotation tables for NCBI submission

Scott M. Geib [1*†], Brian Hall [2†], , Theodore Derego [1], Forest T. Bremer,[2], Kyle Cannoles[1,3], and Sheina B. Sim [1]

[1]Tropical Plant Protection Research Unit, USDA-ARS Daniel K. Inouye U.S. Pacific Basin Agricultural Research Center, Hilo, HI, 96720, USA

[2]Plant and Environmental Protection Science, University of Hawaii at Manoa, Honolulu, HI, 96822, USA.

[3]Department of Computer Science and Engineering, University of Hawaii at Hilo, Hilo, HI, 96720, USA.

*Corresponding author †Authors contributed equally

email: SMG: scott.geib@ars.usda.gov BH:bhall7@hawaii.edu TD:t.derego@yahoo.com FTB:forestb@hawaii.edu KC:kylecann@hawaii.edu SBS:sheina.sim@ars.usda.gov

February 20, 2018

## Abstract

**Background**

One of the most overlooked, yet critical components of a whole genome sequencing project is the submission and curation of the data to a genomic repository, most commonly NCBI. While large genome centers or genome groups have developed software tools for post-annotation assembly filtering, annotation, and conversion into NCBI's annotation table format, these tools typically require back-end setup and connection to an SQL database and/or some knowledge of programming (Perl, Python) to implement. With whole genome sequencing becoming commonplace, genome sequencing projects are moving away from the genome centers, and into the ecology or biology lab, where much less resources are present to support the process of genome assembly curation. To fill this gap, we developed software to assess, filter, transfer annotations, and convert a draft genome assembly and annotation set into NCBI annotation table (.tbl) format, facilitating submission to NCBI Genome Assembly database. This software has no dependencies, is compatible across platforms, and utilizes a simple command to perform a variety of simple and complex post-analysis, pre-NCBI submission WGS project tasks.

**Findings**

The Genome Annotation Generator is a consistent and user-friendly bioinformatics tool that can be used to generate a *.tbl* file that is consistent with the NCBI submission pipeline.

**Conclusions**

The Genome Annotation Generator achieves the goal of providing a publicly available tool that will facilitate the submission of annotated genome assemblies to NCBI. It is useful for any individual researcher or research group who wishes to submit a genome assembly of their study system to NCBI.

# 1   Introduction

While ever-improving sequencing technology and assembly software enable the collection of raw sequences for genome assembly and structural annotation, further steps need to be taken to ensure the quality and completeness of a WGS project for submission to the National Center for Biotechnology Information (NCBI) or other data repositories [36]. To submit a genome to the NCBI for curation, it must be converted to the NCBI annotation table format (*.tbl*). With a genome assembly project consisting of thousands of sequences demarcated by hundreds of thousands of structural annotations, this task clearly requires automation. However, there is currently no freely available tool which performs rapid and controlled conversion of a genome assembly and associated structural annotations into a *.tbl* format in addition to allowing for editing, modification, and revision of the content of the project. Moreover, the typical assembly and draft annotation contains some degree of questionable or erroneous data which requires correction or omission. It may also be desirable to add functional annotations to the submission and integrate results from InterProScan, BLAST homology to curated databases, or ontology terms generated by other tools [20, 5, 22].

The traditional approach used to address these problems is to use Linux command line tools or write custom scripts which modify and filter the genome using a scripting language such as Perl or Python [4, 32, 15] or large scale genomic database systems [25]. This method may not be easily or readily reproducible, or it may be entirely beyond the ability of an investigator who has less familiarity with generating custom scripts *de novo*. Even amongst those researchers who use best practices to write clean, well-tested, and reusable scripts to accomplish these tasks, doing requires a large amount of duplicated effort. For this reason, the Genome Annotation Generator (GAG) was written to provide a straightforward and consistent tool for addressing the most common errors in genome assemblies, adding functional annotations from disparate sources, and producing an NCBI submission-ready annotation *.tbl* file. In addition, the software provides a means for integrating existing functional annotations and marking annotations that require manual curation or review. All of these tasks are done through an intuitive command line program requiring only basic unix skills, and has no required dependencies or packages. The program GAG facilitates the submission of whole genome sequencing (WGS) projects to NCBI as well as provide a standardized utility and workflow that fosters consistency

between projects. Due to emerging genome sequencing initiatives such as the 5,000 Insect Genomes Initiative (i5K), the Plant Genome Initiative, and Genome 10K [19, 26], many independent research groups which are not specialized in genome annotation and analysis are generating large genomic datasets and performing genome sequencing projects within their lab. This program can assist in ensuring quality and consistency of data for new genome biologists.

# 2   Overview

The GAG program is a command line Python program, written in Python 2.7 and requiring no additional outside programs or packages to run. The user directs the program to the genome *.fasta* file and a *.gff3* file containing structure annotations. In addition, a number of options can be used to fix possible errors, flag or remove features (i.e. genomic elements described in the *.gff* structural anntoation file) based on selected criteria, add functional annotations, trim regions of the genome out of the assembly, and, of course, write the genome to NCBI *.tbl* file format. In addition, changes made to the genome annotation, functional annotations added, or flags requesting manual review are also annotated back to the *.gff3* structural annotation file, and the original fasta file is corrected as needed. When the user issues commands to modify the genome, e.g. to remove short introns, the statistics will display two columns, representing the original and modified genomes. This allows for stepwise and documented filtering and review to occur, and interactions between GAG and visual genome review tools (e.g. Artemis, Apollo, GBrowse, JBrowse etc) [35, 29, 21, 31, 28].

# 3   Methods

As an example, we consider a possible work-flow for a user wishing to prepare a genome for submission to the NCBI Eukaryotic WGS Database. The user has a scaffolded genome assembly produced by one of many whole genome assemblers [2, 16, 30] in *.fasta* file format and a corresponding GFF3 feature file [9, 8] containing structural annotations resulting from an automated annotation pipeline or predictors such as Maker, Evidence Modeler, Jigsaw, or others e.g. [3, 18, 1, 17, 33, 34, 7, 10]. The approach would be to first possibly generate functional annotations of predicted genes if this is desired, using whatever approach the user is interested in, and then using the genome and annotations with GAG. After using GAG to remove or flag features of interest, the user then may then further investigate flagged features in a genome browser by loading the output of GAG, edit, and then perform further filtering in GAG, and iterate through this process until a final draft genome product is generated. Finally GAG writes a NCBI table file, on which *tbl2asn* is run for submission to NCBI. This may identify regions of the genome

3

67 that need to be trimmed, due to possible adapter contamination in the genome, or low
68 quality sequence. Any errors generated by *tbl2asn* can then be corrected in GAG, the
69 genome trimmed, until an error free submission is generated.

70     To use GAG, the user creates a folder containing the genome files (or links to them)
71 and runs *gag.py* from the terminal, with the *.fasta* and *.gff3* files. GAG will write a
72 statistics file, containing infomraiton on the number of each feature type, lengths, and
73 other information that may be useful for the submitter. In our experience, automated
74 genome annotation software frequently produces assemblies containing introns as short
75 as 1 base pair long; if any such features are present, GAG can be run to detect them.
76 It is important to note that while NCBI requires short introns to be removed, cutoffs
77 recommended by NCBI may be more stringent that what you want, as they are set to
78 reduce erroneous data being entered into NCBI. For example, prediction of single base
79 introns might not be errors, and represent true data. It is up to the user to dictate what
80 cutoffs they want to define to remove or flag for manual review. To address these short
81 introns, the user simply applies option *-ris* (Remove_Intron_Shorter_Than) with a value
82 of *10*. GAG will discard any mRNA containing an intron shorter than the minimum of
83 ten. A comparison of the genome content before and after removal is printed to the
84 *.stats* file. If the user instead wishes to only flag features that meet these criteria and
85 not remove them, alternatively the *-fis* (Flag_Introns_Shorter_Than) option could be
86 used, which instead adds a *GAG_FLAG* annotation to the attributes column of the *.gff3*
87 file describing the reason for flagging, allowing manual review of flagged features in a
88 genome browser. GAG will automatically update all parent and child features (gene or
89 CDS entries) to reflect removal of mRNA features. A list of available flag or removal
90 options are listed in Table 1.

91     Another review for submission might be that all coding regions be a minimum length.
92 For this example we use 150 base pairs in length, which is suggested by NCBI [11,
93 12]. To add this additional level of filtering, a second option can be used: *-rcs 150*, to
94 Remove_CDS_Shorter_Than 150 bp. When the genome is written to the output folder,
95 GAG will write a file called *genome.removed.gff* containing all the features left out of
96 the final version). It is important to remember that CDS cutoffs at 150 bp will possibly
97 remove some biologically correct amino acids.

98     GAG supports two straightforward correction, or fix tools. If the user's GFF3 file
99 does not explicitly indicate the presence of start and stop codons, or if there is reason to
100 believe there are errors in ORF prediction in the provided GFF file, GAG can add start
101 and stop feature to the GFF file. The user simply issues the command with the option
102 *–fix_start_stop* and these features will be added to the GFF3 file, and their existence
103 noted in the table file. A second issue that can arise in a draft genome assembly is for
104 a contig or scaffold to have a string of ambiguous bases (N's) at the very beginning or
105 end of the contig. These should be removed from the assembly, and can be using the –

106 *fix_terminal_ns* option, as they can be mis-interpreted as scaffold gaps. Removing these

107 regions from the genome though, will disrupt the parity between coordinates in the *.fasta*

108 genome file and the *.gff3* annotation file. GAG will automatically update coordinates

109 in the *.gff3* file to reflect any regions removed from the sequence file. During execution

110 of *tbl2asn* or submission to NCBI, it may be identified that regions of the genome may

111 be contaminated with microbial, vector, or sequencing adapter sequence as part of the

112 "contaminate screen" step. A *.bed* formatted file can be supplied with the *-trim* option,

113 containing regions of the assembly to exclude, either ranges within a contig or scaffold,

114 or an entire scaffold. GAG will update both the *.fasta* and *.gff3* files so that coordinate

115 are still synchronized. This is a particularly difficult operation to perform without a

116 specialized tool.

117 At present, GAG has simple commands to remove or flag introns, exons, coding regions

118 and genes based on minimum or maximum lengths, which will also edit or remove any

119 parent or child feature from the annotation file so as not to create incomplete feature

120 annotations. It can also remove features from a list, which is useful for cases where

121 a genome submission is rejected and a list of invalid mRNAs and genes provided. In

122 addition, all discarded features are retained in a "genome.removed.gff" file and the entire

123 editing session is documented so that the user can retain the filtering criteria used on the

124 particular dataset.

125 GAG supports two methods to add functional annotations to a genome. First, it can

126 read an annotated GFF3 file containing gene names, protein products, cross-references

127 to databases, and ontology terms following GFF3 qualified nomenclature in the *attribute*

128 column of the GFF3 file [24, 23, 27, 6]. Any annotations present will be automatically

129 carried over to the NCBI feature table file. For users with annotations from another

130 source, GAG can read them from a simple tab-delimited file. The annotations supported

131 by the current version of GAG are *Name* (for genes), *Dbxref*, *Ontology_term* and *product*

132 (for descriptive mRNA products). These are also written to a new GFF3 file, so GAG

133 can be utilized as a tool to also functionally annotate a GFF3 file. Detailed instruc-

134 tions for running GAG, examples for each of the main functions (e.g. removing features,

135 adding start and stop codons, trimming features, adding annotations) as well as formats

136 and conversion tools for functional annotations are available on the GAG software web-

137 site webpage: `http://genomeannotation.github.io/GAG/` and a stable release (version

138 2.0.1) is available in an accompanying GigaScience Database entry [13, 14] and is shared

139 through SciCrunch under RRID SCR_016053.


# 4   Implementation

141 GAG is written in Python 2.7. It has no dependencies beyond the standard library. The

142 program is modular, abstracting biological concepts such as Sequence, Gene and CDS

5

143 into classes which may be incorporated into other software tools. In addition, the code
144 is covered by a suite of unit and integration tests, allowing developers to modify or add
145 to the code base with reduced risk of introducing errors. It should be easily executable
146 by the novice programmer with only basic command-line experience, but also powerful
147 enough to be implemented within robust genomic data processing pipelines.

# 5   Conclusion

149 GAG can be easily expanded in the future to support more specific needs of researchers,
150 less common annotation types, and integrate conversion of common functional annotation
151 output formats (e.g. InterProScan, BLAST, Blast2Go) for addition to NCBI annotation
152 table formats. Currently, GAG is an intermediate, but critical tool, between a simple
153 format conversion tool and more sophisticated annotation editors. In future developments
154 of GAG, we plan to allow the integration of multiple lines of evidence supporting gene
155 models to help users discriminate apparently high quality annotations from annotations
156 with little support or possible errors. This could rapidly improve and standardize manual
157 annotation efforts in systems and user groups that are not integrated into genome center
158 annotation pipelines.

# 6   Data and Software

160 Availability and requirements

161   • Project name: Genome Annotation Generator

162   • Project home page: `https://github.com/genomeannotation/GAG`

163   • Operating systems: Linux, Windows, OS

164   • Programming language: Python

165   • Other requirements: Python 2

166   • License: MIT

167   • SciCrunch RRID: SCR_016053

168   Availability of Supporting Data Snapshots of the software and accompanying files are
169 available from the GigaScience GigaDB database, and the algorithm is also demonstrated
170 in Code Ocean repository `https://codeocean.com/2018/02/06/genome-annotation-generator-ncbi-f`
171 [13] [14].

# 7 Declarations

## 7.1 Competing Interests

The authors declare that they have no competing interests

## 7.2 Funding

## 7.3 Authors' contributions

SMG conceived software concept. BH, TD, and SMG designed and wrote software. BH, SMG, and SBS wrote manuscript.

## 7.4 Acknowledgements

7

Table 1: **Options for GAG**

| Option | Type of function | Description |
|---|---|---|
| *-a* <annotation file> | Annotate | Adds functional annotations present in annotation file to *.gff* and *.tbl* |
| *-t* <*.bed* file> | Trim | Removes regions of genome indicated in *.bed* file from *.fasta* and *.gff3* |
| *-fix_start_stop* <no value> | Fix | Adds or corrects start and stop codon features to *.gff3* |
| *-fix_terminal_ns* <no value> | Fix | Removes any trailing ends from contig ends in assembly, updates *.gff3* coordinates |
| *-rcs* <integer> | Remove | Remove CDS shorter than <integer> |
| *-rcl* <integer> | Remove | Remove CDS longer than <integer> |
| *-res* <integer> | Remove | Remove exons shorter than <integer> |
| *-rel* <integer> | Remove | Remove exons longer than <integer> |
| *-ris* <integer> | Remove | Remove introns shorter than <integer> |
| *-ril* <integer> | Remove | Remove introns longer than <integer> |
| *-rgs* <integer> | Remove | Remove genes shorter than <integer> |
| *-rgl* <integer> | Remove | Remove genes longer than <integer> |
| *-fcs* <integer> | Flag | Flag CDS shorter than <integer> |
| *-fcl* <integer> | Flag | Flag CDS longer than <integer> |
| *-fes* <integer> | Flag | Flag exons shorter than <integer> |
| *-fel* <integer> | Flag | Flag exons longer than <integer> |
| *-fis* <integer> | Flag | Flag introns shorter than <integer> |
| *-fil* <integer> | Flag | Flag introns longer than <integer> |
| *-fgs* <integer> | Flag | Flag genes shorter than <integer> |
| *-fgl* <integer> | Flag | Flag genes longer than <integer> |

# References

[1] Jonathan E. Allen and Steven L. Salzberg. Jigsaw: integration of multiple sources of evidence for gene prediction. *Bioinformatics*, 21(18):3596–3603, 2005.

[2] J. Butler, I. MacCallum, M. Kleber, I. A. Shlyakhter, M. K. Belmonte, E. S. Lander, C. Nusbaum, and D. B. Jaffe. Allpaths: de novo assembly of whole-genome shotgun microreads. *Genome Res*, 18(5):810–20, 2008.

[3] Brandi L. Cantarel, Ian Korf, Sofia M. C. Robb, Genis Parra, Eric Ross, Barry Moore, Carson Holt, Alejandro Sanchez Alvarado, and Mark Yandell. Maker: An easy-to-use annotation pipeline designed for emerging model organism genomes. *Genome Research*, 18(1):188–196, 2008.

[4] Peter J. A. Cock, Tiago Antao, Jeffrey T. Chang, Brad A. Chapman, Cymon J. Cox, Andrew Dalke, Iddo Friedberg, Thomas Hamelryck, Frank Kauff, Bartek Wilczynski, and Michiel J. L. de Hoon. Biopython: freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):1422–1423, 2009.

[5] Ana Conesa, Stefan Götz, Juan Miguel García-Gómez, Javier Terol, Manuel Talón, and Montserrat Robles. Blast2go: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics*, 21(18):3674–3676, 2005.

[6] The Gene Ontology Consortium. Expansion of the gene ontology knowledgebase and resources. *Nucleic Acids Research*, 45(D1):D331–D338, 2017.

[7] Val Curwen, Eduardo Eyras, T. Daniel Andrews, Laura Clarke, Emmanuel Mongin, Steven M.J. Searle, and Michele Clamp. The ensembl automatic gene annotation system. *Genome Research*, 14(5):942–950, 2004.

[8] K. Eilbeck and S. E. Lewis. Sequence ontology annotation guide. *Comp Funct Genomics*, 5(8):642–7, 2004.

[9] Karen Eilbeck, Suzanna Lewis, Christopher Mungall, Mark Yandell, Lincoln Stein, Richard Durbin, and Michael Ashburner. The sequence ontology: a tool for the unification of genome annotations. *Genome Biology*, 6(5):R44, 2005.

[10] Christine G. Elsik, Kim C. Worley, Anna K. Bennett, Martin Beye, Francisco Camara, Christopher P. Childers, Dirk C. de Graaf, Griet Debyser, Jixin Deng, Bart Devreese, Eran Elhaik, Jay D. Evans, Leonard J. Foster, Dan Graur, Roderic Guigo, Katharina Jasmin Hoff, Michael E. Holder, Matthew E. Hudson, Greg J. Hunt, Huaiyang Jiang, Vandita Joshi, Radhika S. Khetani, Peter Kosarev, Christie L.

Kovar, Jian Ma, Ryszard Maleszka, Robin F. A. Moritz, Monica C. Munoz-Torres, Terence D. Murphy, Donna M. Muzny, Irene F. Newsham, Justin T. Reese, Hugh M. Robertson, Gene E. Robinson, Olav Rueppell, Victor Solovyev, Mario Stanke, Eckart Stolle, Jennifer M. Tsuruda, Matthias Van Vaerenbergh, Robert M. Waterhouse, Daniel B. Weaver, Charles W. Whitfield, Yuanqing Wu, Evgeny M. Zdobnov, Lan Zhang, Dianhui Zhu, and Richard A. Gibbs. Finding the missing honey bee genes: lessons learned from a genome upgrade. *BMC Genomics*, 15(1):86, 2014.

[11] National Center for Biotechnology Information. The genbank submissions handbook [internet], 2011.

[12] National Center for Biotechnology Information. Common discrepancy reports, January 2013.

[13] S. M. Geib, B. Hall, T. Derego, F. T. Bremer, K. Cannoles, and S. B. Sim. Supporting data for "genome annotation generator: A simple tool for generating and correcting wgs annotation tables for ncbi submission". *GigaScience Database*, 2017. `http://dx.doi.org/10.5524/100308`.

[14] S. M. Geib, B. Hall, T. Derego, F. T. Bremer, K. Cannoles, and S. B. Sim. Genome annotation generator ncbi for submission [source code]. *CodeOcean*, 2018. `https://doi.org/10.24433/CO.fceb0521-a26d-441f-9fe0-bccc6a250fc9`.

[15] Robert Gentleman, Vincent Carey, Douglas Bates, Ben Bolstad, Marcel Dettling, Sandrine Dudoit, Byron Ellis, Laurent Gautier, Yongchao Ge, Jeff Gentry, Kurt Hornik, Torsten Hothorn, Wolfgang Huber, Stefano Iacus, Rafael Irizarry, Friedrich Leisch, Cheng Li, Martin Maechler, Anthony Rossini, Gunther Sawitzki, Colin Smith, Gordon Smyth, Luke Tierney, Jean Yang, and Jianhua Zhang. Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology*, 5(10):R80, 2004.

[16] Sante Gnerre, Iain MacCallum, Dariusz Przybylski, Filipe J. Ribeiro, Joshua N. Burton, Bruce J. Walker, Ted Sharpe, Giles Hall, Terrance P. Shea, Sean Sykes, Aaron M. Berlin, Daniel Aird, Maura Costello, Riza Daza, Louise Williams, Robert Nicol, Andreas Gnirke, Chad Nusbaum, Eric S. Lander, and David B. Jaffe. High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proceedings of the National Academy of Sciences*, 108(4):1513–1518, 2011.

[17] Brian Haas, Steven Salzberg, Wei Zhu, Mihaela Pertea, Jonathan Allen, Joshua Orvis, Owen White, C Robin Buell, and Jennifer Wortman. Automated eukaryotic gene structure annotation using evidencemodeler and the program to assemble spliced alignments. *Genome Biology*, 9(1):R7, 2008.

[18] Carson Holt and Mark Yandell. Maker2: an annotation pipeline and genome-database management tool for second-generation genome projects. *Bmc Bioinformatics*, 12, 2011.

[19] i5K Consortium. The i5k initiative: Advancing arthropod genomics for knowledge, human health, agriculture, and the environment. *Journal of Heredity*, 104(5):595–600, 2013.

[20] Philip Jones, David Binns, Hsin-Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla, Hamish McWilliam, John Maslen, Alex Mitchell, Gift Nuka, Sebastien Pesseat, Antony F. Quinn, Amaia Sangrador-Vegas, Maxim Scheremetjew, Siew-Yit Yong, Rodrigo Lopez, and Sarah Hunter. Interproscan 5: genome-scale protein function classification. *Bioinformatics*, 30(9):1236–1240, 2014.

[21] Eduardo Lee, Gregg A. Helt, Justin T. Reese, Monica C. Munoz-Torres, Chris P. Childers, Robert M. Buels, Lincoln Stein, Ian H. Holmes, Christine G. Elsik, and Suzanna E. Lewis. Web apollo: a web-based genomic annotation editing platform. *Genome Biology*, 14(8):R93, 2013.

[22] Michele Magrane and UniProt Consortium. Uniprot knowledgebase: a hub of integrated protein data. *Database*, 2011, 2011.

[23] Barry Moore, Guozhen Fan, and Karen Eilbeck. Soba: sequence ontology bioinformatics analysis. *Nucleic Acids Research*, 38(suppl 2):W161–W164, 2010.

[24] Christopher J. Mungall, Colin Batchelor, and Karen Eilbeck. Evolution of the sequence ontology terms and relationships. *Journal of Biomedical Informatics*, 44(1):87–93, 2011.

[25] Christopher J. Mungall, David B. Emmert, and The FlyBase Consortium. A chado case study: an ontology-based modular schema for representing genome-associated biological information. *Bioinformatics*, 23(13):i337–i346, 2007.

[26] Genome 10K Community of Scientists. Genome 10k: A proposal to obtain whole-genome sequence for 10,000 vertebrate species. *Journal of Heredity*, 100(6):659–674, 2009.

[27] Martin Reese, Barry Moore, Colin Batchelor, Fidel Salas, Fiona Cunningham, Gabor Marth, Lincoln Stein, Paul Flicek, Mark Yandell, and Karen Eilbeck. A standard variation file format for human genome sequences. *Genome Biology*, 11(8):R88, 2010.

[28] J. T. Robinson, H. Thorvaldsdottir, W. Winckler, M. Guttman, E. S. Lander, G. Getz, and J. P. Mesirov. Integrative genomics viewer. *Nat Biotechnol*, 29(1):24–6, 2011.

[29] K. Rutherford, J. Parkhill, J. Crook, T. Horsnell, P. Rice, M. A. Rajandream, and B. Barrell. Artemis: sequence visualization and annotation. *Bioinformatics*, 16(10):944–5, 2000.

[30] Jared T. Simpson, Kim Wong, Shaun D. Jackman, Jacqueline E. Schein, Steven J.M. Jones, and Inanc Birol. Abyss: A parallel assembler for short read sequence data. *Genome Research*, 19(6):1117–1123, 2009.

[31] Mitchell E. Skinner, Andrew V. Uzilov, Lincoln D. Stein, Christopher J. Mungall, and Ian H. Holmes. Jbrowse: A next-generation genome browser. *Genome Research*, 19(9):1630–1638, 2009.

[32] Jason E. Stajich, David Block, Kris Boulez, Steven E. Brenner, Stephen A. Chervitz, Chris Dagdigian, Georg Fuellen, James G.R. Gilbert, Ian Korf, Hilmar Lapp, Heikki Lehväslaiho, Chad Matsalla, Chris J. Mungall, Brian I. Osborne, Matthew R. Pocock, Peter Schattner, Martin Senger, Lincoln D. Stein, Elia Stupka, Mark D. Wilkinson, and Ewan Birney. The bioperl toolkit: Perl modules for the life sciences. *Genome Research*, 12(10):1611–1618, 2002.

[33] Mario Stanke, Oliver Schoffmann, Burkhard Morgenstern, and Stephan Waack. Gene prediction in eukaryotes with a generalized hidden markov model that uses hints from external sources. *BMC Bioinformatics*, 7(1):62, 2006.

[34] Mario Stanke and Stephan Waack. Gene prediction with a hidden markov model and a new intron submodel. *Bioinformatics*, 19(suppl 2):ii215–ii225, 2003.

[35] Lincoln D. Stein, Christopher Mungall, ShengQiang Shu, Michael Caudy, Marco Mangone, Allen Day, Elizabeth Nickerson, Jason E. Stajich, Todd W. Harris, Adrian Arva, and Suzanna Lewis. The generic genome browser: A building block for a model organism system database. *Genome Research*, 12(10):1599–1610, 2002.

[36] Mark Yandell and Daniel Ence. A beginner's guide to eukaryotic genome annotation. *Nature Reviews Genetics*, 13(5):329–342, 2012.